

- Since probabilities depend on weights (minimizing weighted loss) and weights depend on probabilities we use
 - iterative weighted least squares / Fisher scoring
 to decide weights/probs (ends up being the same as gradient descent)

Multidimensional Probability:

Have random vector $\vec{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix}$, $\Pr(\vec{X} \in A) = \int_A f(\vec{x}) d\vec{x} = \int_A f(x_1, \dots, x_p) dx_1 \dots dx_p$

- probability that \vec{X} is within $\pm h$ of $\vec{x} \approx h^p f(\vec{x})$

$$E[\vec{X}] = \int \vec{x} f(\vec{x}) d\vec{x} = \begin{bmatrix} E[x_1] \\ \vdots \\ E[x_p] \end{bmatrix}, \quad \text{Var}[\vec{X}] = \begin{bmatrix} \text{Var}(x_1) & \text{Cov}(x_1, x_2) & \dots & \text{Cov}(x_1, x_p) \\ \text{Cov}(x_2, x_1) & \text{Var}(x_2) & \dots & \text{Cov}(x_2, x_p) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(x_p, x_1) & \dots & \dots & \text{Var}(x_p) \end{bmatrix}_{p \times p}$$

$\text{Var}[\vec{X}]$ is positive semidefinite: $\vec{u}^T \text{Var}(\vec{X}) \vec{u} \geq 0 \quad \forall \vec{u}$
 and it is symmetric, ~~non null~~

These ~~non null~~ indicate it has an eigendecomposition:

$\text{Var}[\vec{X}] = w d w^T$ where d = diagonal matrix of eigenvalues of $\text{Var}[\vec{X}]$ (all eigenvals are ≥ 0)

w = $p \times p$ matrix whose columns are eigenvectors of $\text{Var}[\vec{X}]$, scaled so that each column has length 1. All eigenvectors are orthogonal to each other so $w^T = w^{-1}$

For vector \vec{a} , $\text{Var}[\vec{a}\vec{X} + \vec{b}] = \vec{a} \text{Var}[\vec{X}] \vec{a}^T$

$$E[d^{-1/2} w^T (\vec{X} - E[\vec{X}])] = \vec{0}$$

$$\text{Var}[d^{-1/2} w^T (\vec{X} - E[\vec{X}])] = I$$

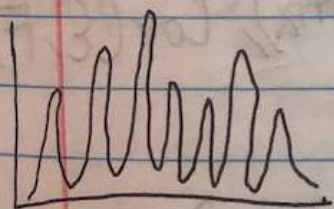
- cond independence
- nonparametric model needs to have more args when more data points are added

Empirical probabilities:

$$\hat{P}(X \in A) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \in A\}$$

Kernel Density Estimation: (estimating probability densities for data features)

$$\hat{f}_{\text{kernel}}(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} k\left(\frac{x - x_i}{h}\right)$$



- large h means less peaks that are wider
- lower h means more and higher peaks

- Use cross validation to pick bandwidth with either $\int (f(x) - \hat{f}(x))^2 dx$ or $\int -\log(\hat{f}(x)) f(x) dx$
integrated squared error log likelihood

(with IID samples for large n , by central limit theorem)

$$\int -\log(\hat{f}(x)) f(x) dx \approx \frac{1}{n} \sum_{i=1}^n -\log \hat{f}(x_i)$$

- If we can find some X_j such that X_1, \dots, X_p are all conditionally independent given X_j , then $f(x_1, \dots, x_p) = f(x_j) \prod_{i=1, i \neq j}^p f(x_i | x_j)$. This reduces number

of param we need to estimate (and helps lessen curse of dimensionality).

Factor Model:

- Assume $E[\vec{X}] = 0$
- Assume there is a q -dimensional random vector \vec{F} with $E[\vec{F}] = 0$, $\text{Var}[\vec{F}] = I$, $q < p$
- \vec{F} is latent/unobserved/hidden, \vec{X} is manifest/observable
- Assume $\vec{X} = W\vec{F} + \vec{\epsilon}$ with $E[\vec{\epsilon}] = 0$, $\text{Var}(\vec{\epsilon}) = \Psi$
($\Psi = \text{diag matrix}$), $\text{Cov}(\vec{\epsilon}, \vec{F}) = 0$

So under these assumptions

$$\text{Cov}(X_i, X_j) = \text{Cov}\left(\sum_{k=1}^q w_{ik} F_k + \epsilon_i, \sum_{l=1}^q w_{jl} F_l + \epsilon_j\right)$$

$$= \text{Cov}\left(\sum_{k=1}^q w_{ik} F_k, \sum_{l=1}^q w_{jl} F_l\right)$$

$$= \sum_{k=1}^q \sum_{l=1}^q w_{ik} w_{jl} \text{Cov}(F_k, F_l)$$

$$= \sum_{k=1}^q w_{ik} w_{jk} \quad (\text{since } \text{Cov}(F_k, F_l) \text{ is } 0 \text{ when } k \neq l \text{ and } 1 \text{ when } k = l)$$

So X_i and X_j are correlated when they load on the same factors